

# A Comparative Analysis of Machine Learning for Restaurant Rating Prediction

Name: Jaat Devrajsingh Bharat Singh  
Dept: Data Science, Manipal University Jaipur  
New Delhi, India  
vs601710@gmail.com

Name: Pramod Prajapati  
Dept: Department of Computer science and Engineering  
prasunetcompany@gmail.com

**Abstract:** This paper offers a thorough investigation into machine learning methods for restaurant rating prediction. As the use of online reviews grows, automated techniques for rating prediction can offer useful information to consumers and company owners alike. To create predictive models, we examine restaurant metadata and customer reviews. In-depth feature engineering, thorough data preprocessing, and the application of several machine learning models such as the Support Vector Classifier (SVC), light bgm, Extra Tree regressor, Random Forest, Bagging Classifier, etc. Are all examined in this study. Accuracy scores, error-based metrics, and statistical correlation metrics are used to assess model performance. We also include comparisons of model performance and visualizations of important data patterns. The findings demonstrate the efficacy of various modeling techniques and offer information on the key characteristics influencing rating predictions.

**Keywords:** Restaurant Ratings, Sentiment Analysis, Machine Learning, Predictive Modeling, Feature Engineering

## I. INTRODUCTION

Online restaurant reviews are a great source of knowledge since they provide thorough understanding of consumer experience and satisfaction. These evaluations are very important for guiding prospective patrons toward wise decisions prior to dining at a restaurant and for giving owners useful comments to enhance their offerings. Restaurant ratings are manually assigned historically, usually based on direct client comments. But as technology develops, automated systems to forecast ratings depending on textual reviews become increasingly necessary. Forecasting restaurant ratings from online reviews is a difficult task requiring several approaches, especially in the domains of natural language processing (NLP) and machine learning. These disciplines seek to unravel the intricate interactions between consumer satisfaction and textual materials.

Accurate prediction of restaurant ratings using machine learning models could simplify and automate the decision-making process for owners of restaurants as well as consumers. These predictive models can provide fast, real-time insights by analyzing vast sets of reviews, so removing the need for human assessment. By spotting trends or repeating patterns in consumer comments, automated rating prediction can also help companies to pinpoint areas of development and enable them to adapt and enhance their offerings more successfully.

Predicting ratings depending just on reviews is not an easy chore, though. It calls for a strong awareness of the sentiment the book conveys, which might be either positive or negative and with many subtleties in between. Thus, this process depends critically on sentiment analysis that is, the classification of word emotional tone. Although conventional rating systems only consider numerical values, sentiment analysis offers a more complex, rich view of consumer experiences, so improving prediction accuracy.

In order to identify the best method for predicting restaurant ratings from online reviews, this study assesses and contrasts various machine learning models. The research aims to offer important insights into how emotions in reviews affect ratings by incorporating sentiment analysis. By comparing predictive performances empirically and making suggestions for improved restaurant decision-making tools, the study advances the field.

## II. DATASET & FEATURES

Each of the 19,896 entries in the dataset in question represents a distinct Yelp platform review. Each entry in this index, which ranges from 0 to 19,895, includes four columns that highlight important features of the user-posted reviews. The URL for each individual review is contained in the first column, Yelp URL, an object-type data field that essentially acts as a link to that particular review on the Yelp platform. The user's overall rating for the company or service under review is represented by the second column, Rating, which is an integer-type field with values ranging from 1 to 5. Thus, this column offers a numerical depiction of customer satisfaction and sentiment.

The precise date and time the review was posted are recorded in the third column, Date, which uses the datetime64 format. For time-based analysis, like looking at seasonality or changes in user sentiment, this column is crucial for comprehending patterns in review activity over time. The written text of the review itself is contained in the last column, Review Text, which is an object-type field. In addition to providing more thorough and contextualized information about the product or service under review, this textual data offers insightful qualitative information about the user experience.

The entire dataset appears to be clean and complete, with each review entry fully filled in all fields, as there are no missing or null values in any of its columns. Furthermore, the dataset's 19,896 rows provide a sizable foundation for statistical and text-based analysis, especially when it comes

to finding trends and connections among the ratings, review text, and submission dates. With a memory usage of roughly 621.9 KB, the dataset is manageable in terms of processing and storage requirements, allowing for extensive analysis using a variety of computational techniques, including correlation analysis to look into possible relationships between the ratings and the reviews' textual content, time series analysis to understand trends, and natural language processing (NLP) for sentiment analysis.

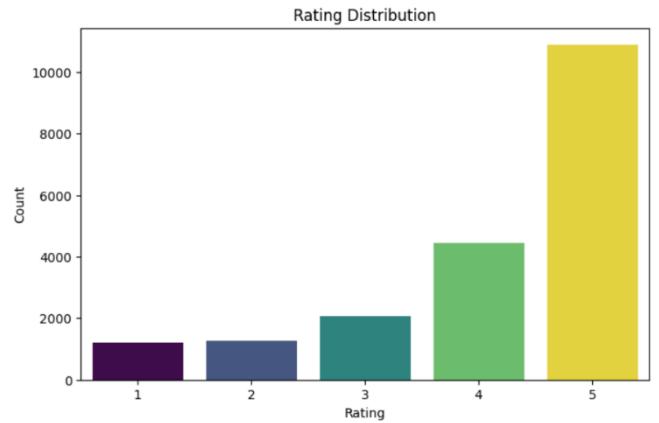
### III. METHODOLOGY

A number of preprocessing procedures were used to clean, transform, and analyze the text data in order to get the Yelp review dataset ready for analysis. By following these procedures, the data was guaranteed to be in a consistent format appropriate for feature extraction, model training, and additional analysis. To better comprehend the format and content of the reviews, exploratory data analysis (EDA) methods were also used.

An initial evaluation of the dataset was conducted to look for duplicate entries, missing values, and inconsistencies. Notably, we were able to move straight forward with exploratory analysis and text preprocessing since no missing or duplicate values were discovered. Textual restaurant reviews with numerical ratings between 1 and 5 make up the majority of the dataset. We analyzed the rating distribution, word frequency, and review length in order to gain a better understanding of the data's distribution and characteristics. This analysis yielded important insights into customer sentiment and review trends. These analyses are described in detail in the sections that follow.

### 1. Data Preprocessing:

**1.1 Rating Distribution:** A frequency distribution plot was created in order to obtain an understanding of the dataset's rating distribution (Figure 1). Given that most reviews have a 5-star rating, the results show a notable class imbalance. While ratings of three and four stars show a gradual increase but are still relatively uncommon, lower ratings, like one or two stars, are much less common. This disparity implies that the dataset is biased in favor of satisfied customers, which could skew predictive modeling. When developing predictive models, it might be necessary to address this class imbalance using methods like weighted loss functions or resampling.



- 1.2 **Word Cloud:** To get a grasp on the most commonly used words in the dataset, we created a word cloud (see Figure 2). This visual representation brings to light key phrases like “ice cream,” “delicious,” “cake,” and “sweet,” hinting that a significant number of reviews focus on desserts and pastries. The abundance of words tied to positive experiences shows that most customers are quite happy with their dining outings. By recognizing these frequently mentioned terms, we can better understand which aspects of a restaurant are often highlighted and what might resonate with customers.

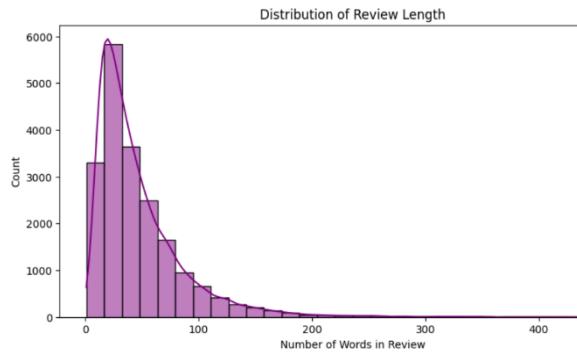
Moreover, the inclusion of words such as “order,” “serve,” and “location” points to the importance of service quality and convenience in shaping customer satisfaction. This understanding can guide us in refining our feature engineering methods for predictive modeling, like using topic modeling or sentiment-driven features.



- 1.3 Review Length Distribution:** To dive into the verbosity of customer reviews, we took a closer look at how review lengths vary, counting the number of words in each review (see Figure 3). The histogram reveals a right-skewed distribution, which means that most reviews tend to be on the shorter side, with a noticeable drop in frequency as the length increases. While we see plenty of brief reviews, the existence of longer ones indicates that some customers are willing to share more detailed feedback, which can provide richer sentiment-based insights.

This observation is critical for text preprocessing, as excessively short reviews may lack sufficient context for sentiment analysis or rating prediction. Conversely, very long reviews could introduce noise.

By setting appropriate length thresholds, we can balance data richness and relevance while ensuring computational efficiency in model training.



## 2. Text preprocessing:

**2.1 Text Lowercasing:** To keep things consistent, we converted all review text to lowercase. This way, we eliminated any confusion that might arise from case sensitivity. For example, "Good" and "good" are now treated as the same word, which helps us avoid discrepancies in our analysis.

**2.2 Punctuation and Noise Removal:** We stripped away punctuation marks (like periods, commas, and exclamation points), special symbols (such as emojis), and numbers from the review text. These elements don't really add any value to sentiment analysis, so we removed them to let the model focus on the words that truly express sentiment. For instance, the review "I love the food!! 🍔🍟" was simplified to "i love the food."

**2.3 Stop Word Removal:** Using the Natural Language Toolkit (NLTK) library, we got rid of stop words like "and," "the," "is," and "to" from the reviews. These common words don't contribute much to sentiment analysis, so removing them helps reduce noise and allows the model to concentrate on more significant terms.

**2.4 Creation of Cleaned\_Text Column:** Once the text was cleaned up, we stored the modified version in a new column called `cleaned_text`. This column holds the review text in its polished form, making it easier for further processing without messing with the original data.

**2.5 Tokenization:** We broke down the review text into individual words or tokens through a process called tokenization. This step is crucial for creating a structure that's easier to analyze and interpret. For example, the sentence "The food was amazing!" gets tokenized into ["the", "food", "was", "amazing"]

**2.6 Lemmatization:** We lemmatized the tokens using the WordNet Lemmatizer from NLTK, which reduces each word to its base form. This means that words like "running," "runner," and "ran" are all treated as the same word, enhancing consistency and cutting down the total number of unique words in our dataset.

**2.7 Extra Space Removal:** We made sure to tidy up the text by removing any extra spaces between words, as well as any leading or trailing spaces. This little cleanup step helped ensure that the text was neat and ready for the next round of analysis.

## 3. Feature Engineering:

Feature engineering is essential for turning raw data into a format that machine learning models can work with. In this study, we applied feature engineering to the Yelp reviews dataset through several important steps: Sentiment Analysis, TF-IDF Transformation, and Correlation Analysis. Each of these steps helped us pull out valuable insights from the text reviews, which we then used to predict restaurant ratings.

**3.1 Sentiment Analysis:** We kicked things off with sentiment analysis on the review text to extract two key features: Polarity and Subjectivity. For this, we used the TextBlob library, a well-known tool in the natural language processing (NLP) world that makes it easy to perform sentiment analysis on text.

- **Polarity:** Polarity is all about measuring the sentiment in the review. It ranges from -1 to +1, where -1 means a negative sentiment, 0 is neutral, and +1 indicates a positive sentiment. Essentially, polarity tells us how favorable or unfavorable the sentiment is towards the business or service being reviewed.
- **Subjectivity:** On the other hand, Subjectivity gauges how opinionated or subjective the review is. It ranges from 0 to 1, with 0 being an objective statement (factual or neutral) and 1 being a highly subjective or opinionated statement. Reviews that score higher on subjectivity often include personal opinions, feelings, or experiences, while those with lower subjectivity values tend to stick to more factual descriptions.

**3.2 TF-IDF Transformation:** To turn the text data into a numerical format that machine learning models can work with, we applied a TF-IDF (Term Frequency-Inverse Document Frequency) transformation to the `cleaned_text` column. TF-IDF is a method that highlights the significance of a word within the context of the entire collection of documents (in this case, reviews). We used the TF-IDF vectorizer to convert the text data into a numerical format, assigning a weight to each word based on how often it appears in a specific review and how uncommon it is across all the reviews.

*The key steps involved in the TF-IDF transformation are:*

- **Term Frequency (TF):** This measures how often a term shows up in a particular document (review).

The more a word appears in a document, the higher its TF score will be.

- **Inverse Document Frequency (IDF):** This measures how rare or common a term is across the whole collection. If a word is frequently found in all reviews (like "the" or "and"), it gets a lower IDF score, while less common words that show up in fewer reviews earn a higher IDF score.

The **TF-IDF formula** is expressed as:

$$TF - IDF(t, d) = TF(t, d) \times \log\left(\frac{N}{DF(t)}\right)$$

Where:

t is the term (word),

d is the document (review),

N is the total number of documents (reviews),

DF(t) is the document frequency of the term t.

This transformation turns the `cleaned_text` column into a numerical matrix with 5000 features (terms), where each row represents a review, and each column indicates the TF-IDF score for a specific term in the vocabulary.

- The **TF-IDF** matrix was then combined with the **polarity** and **subjectivity** features to create the final feature matrix:

$$X = TF - IDF \text{ feature} \cup \text{Polarity and Subjectivity}$$

- The target variable, **rating**, was stored in **y**, and this formed the dataset used for model training.

**3.3 Correlation Analysis:** Before we dive into applying the models, we conducted a correlation analysis to explore the relationships between the sentiment features (polarity and subjectivity) and the Rating. We calculated Pearson's Correlation and Spearman's Rank Correlation to get a sense of the strength and nature of these relationships:

**Pearson Correlation:** 0.5205

**Spearman Correlation:** 0.4996

The Pearson correlation reveals a moderate positive linear relationship between polarity and rating, indicating that as sentiment polarity becomes more positive, the rating tends to rise as well. Meanwhile, the Spearman correlation reflects a similar trend but focuses on rank-based (non-linear) relationships, which further reinforces this connection.

#### 4. Machine Learning Model

In this study, we explored various machine learning models to predict restaurant ratings using engineered features like

sentiment polarity, subjectivity, and TF-IDF values. The upcoming sections will dive into each model, detailing how they work and then evaluating their performance metrics.

**4.1 Linear Regression:** Linear Regression is a fundamental machine learning technique used to predict continuous values, which in this context means restaurant ratings. It operates on the assumption that there's a linear relationship between the target variable and the input features. The basic formula for linear regression is:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

Where:

- y is the dependent variable (rating),
- $\beta_0$  is the intercept (constant),
- $\beta_1, \beta_2, \dots, \beta_n$  are the regression coefficients (weights for each feature),
- $x_1, x_2, \dots, x_n$  are the input features (polarity, subjectivity, and TF-IDF),
- $\epsilon$  is the error term.
- 

Linear regression finds the optimal coefficients  $\beta$  by minimizing the **mean squared error** (MSE), which is given by:

$$MSE = \frac{1}{N} = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Where  $y_i$  is the actual value, and  $\hat{y}_i$  is the predicted value.

**Performance:**

- **Root Mean Square Error (RMSE):** 1.0372
- **R<sup>2</sup> (Coefficient of Determination):** 0.2656

**Performance Summary:** The model explained only about 26.56% of the variance in the ratings (low R<sup>2</sup>), suggesting that linear regression was not fully effective at capturing the complexity of the data.

**4.2 Random Forest:** Random Forest is an ensemble method that builds multiple decision trees during the training phase and combines their results to enhance accuracy and reduce the risk of overfitting. Each tree is trained on a random subset of the training data (a process known as bootstrapping), and at each split in the tree, a random selection of features is considered. This element of randomness helps ensure that the trees are varied and minimizes the chances of overfitting.

The Random Forest algorithm is based on the following key steps:

- **Bootstrap Sampling:** Generate several different training datasets by sampling with replacement from the original dataset.

- **Feature Randomization:** At each decision tree split, select a random subset of features to use for the split.
- **Tree Building:** Construct multiple decision trees using these randomized datasets.
- **Voting:** For classification tasks, predictions are determined by majority voting; for regression tasks, the predictions are averaged across all trees.
- Random Forest effectively reduces variance and prevents overfitting by aggregating the predictions from multiple decision trees. The final prediction is calculated as the average (for regression tasks) of the outputs from all individual trees.

**Performance:**

**Accuracy (Before Hyperparameter Tuning):** 59.17%

**Accuracy (After Grid Search):** 56.23%

**Performance Summary:** While Random Forest showed a solid baseline performance with an accuracy of 59.17%, tuning the hyperparameters resulted in a drop in performance, indicating possible overfitting or less-than-ideal parameter choices during the tuning process.

### 4.3 XGBOOST (EXTREME GRADIENT BOOSTING)

Extreme Gradient Boosting, or XGBoost for short, is a robust ensemble machine learning technique that's all about speed and performance. It uses a gradient boosting framework to build a series of decision trees one after the other, where each new tree aims to fix the mistakes made by the previous ones. XGBoost also features advanced regularization techniques (L1 and L2), parallel computation, and built-in tools for dealing with missing values, making it an excellent option for high-dimensional and structured data, including features derived from text.

In this study, we applied XGBoost to predict restaurant ratings by combining sentiment polarity, subjectivity, and TF-IDF features pulled from customer reviews. Given the class imbalance in our dataset, we explored three different versions of XGBoost: the standard model, one that used computed class weights, and another that was trained on a SMOTE-balanced dataset.

#### 4.3.1 XGBoost (Default Configuration)

The baseline XGBoost model was trained without any tweaks for class imbalance. It utilized the complete set of engineered features and operated with default settings, requiring only minimal adjustments.

**Performance:**

- Accuracy: 64.00%

**Summary:** The default configuration of XGBoost delivered the best accuracy among the three variants we tested, achieving an impressive 64.00%. This solid baseline performance highlights XGBoost's capability to effectively model non-linear relationships and identify complex patterns in high-dimensional feature spaces, even when faced with moderate class imbalance. Its boosting framework seems to generalize well across the rating distribution, making it a dependable choice for this task.

#### 4.3.2 XGBoost with compute class weight

To tackle the uneven distribution of rating classes, we calculated and assigned class weights that were inversely proportional to how often each class appeared. This means that the model would face a bigger penalty for getting minority class instances wrong, which encourages it to treat all classes more fairly.

Even though XGBoost doesn't directly support the class\_weight parameter, we managed to balance the classes by tweaking instance weights during training according to the class distribution.

**Performance:**

**Accuracy:** 60.80%

**Summary:**

By applying the calculated class weights, we saw a drop in accuracy to 60.80%, compared to the default setup. This indicates that while the model became more attuned to the less represented classes, it may have compromised its performance on the majority classes, resulting in an overall decrease in predictive accuracy. This outcome points to a possible trade-off between being sensitive to class imbalances and maintaining general performance, especially when weighting isn't combined with other balancing methods.

#### 4.3.3 XGBoost with SMOTE

To further combat class imbalance, we used SMOTE to oversample the minority classes by creating synthetic instances based on the similarities in the feature space of existing samples. This approach helped artificially balance the dataset before training the model, giving it more chances to learn from the underrepresented classes.

**Performance:**

- **Accuracy:** 62.46%

**Summary:**

Implementing SMOTE resulted in a significant improvement over the class-weighted model, raising accuracy to 62.46%. However, it still didn't quite reach the default model's accuracy of 64.00%. This suggests that while SMOTE effectively enhanced representation across classes, the synthetic examples might not have fully captured the complexity of real reviews. Still, the increase from the class-weighted model shows that data-level balancing techniques can be quite beneficial for dealing with class imbalances.

#### 4.4 Extra Tree Classifier

The Extra Trees Classifier (Extremely Randomized Trees) is an ensemble learning technique that is implemented to increase the accuracy of a model while decreasing overfitting. The Extra Trees Classifier is very similar to the Random Forest algorithm and both algorithms use decision trees to complete classification tasks. However, the Extra Trees classifier incorporates more randomness at the point of splitting nodes in each decision tree, compared to the process used in Random Forests (which is more deterministic in nature).

In a simple decision tree, the algorithm will select the best split based on a feature and threshold to maximize some criteria (i.e., Gini impurity or Information Gain). The Extra Trees Classifier takes this novelty to another level, since it first selects a random subset of features and then splits the features on random thresholds before evaluating. This variate gives us trees that are less correlated with each other to reduce the variance of the model and provide more generalizable models.

The Extra Trees Classifier has two main randomization methods:

1. **Random Feature Selection:** When creating trees, the algorithm randomly chooses a set of features to use for each split.
2. **Random Threshold Selection:** For every feature that the algorithm selects, it randomly selects a threshold, rather than just using the best threshold from the data splits.

By adding even more randomness, the model can be less prone to overfitting/reporting trends that the model will not see in the future - this is especially true when the training data includes features that are noisy or less informative. In this sense, we end up with a more robust model that should perform well on unseen data.

#### Performance Evaluation

The Extra Trees Classifier was tested on a [whatever the data set or task was], which had the following performance metric:

- **Accuracy:** 59.77%

Accuracy as a performance metric is expressed as the ratio of correct predictions to total predictions. While the accuracy is a good indication of the classifier's effectiveness at solving the problem, given the class imbalance present in the original dataset, it is often important to look at other performance metric (to be format examples: precision, recall, and f1) to better understand its predictive performance.

#### 4.5 Light GBM

The Light Gradient Boosting Machine (LightGBM) is a well-performing gradient boosting framework developed to increase the training speed and scalability of traditional gradient boosting models. The LightGBM approach varies from conventional gradient boosting in that it grows trees leaf-wise rather than level-wise. This method of growth can provide a more accurate model with fewer trees. Compared with traditional gradient boosting, LightGBM uses a histogram-based learning algorithm, which effectively allows it to work with larger datasets by affecting continuous feature values with histograms.

Traditional gradient boosting builds each tree in order, where each new tree corrects overarching issues from previous trees. The LightGBM tree growth method grows the tree by picking the leaf that leads to the largest reduction in overall loss from the previous trees. This is different from traditional gradient boosting, or from algorithms like XGBoost, which use level-wise growth meaning that the trees split at each level for each leaf in the model. This tree growth method may provide better purity and therefore more accuracy but it can also lead to overfits due to the growth method from which the errors are propagated if the dataset is smaller.

Furthermore, LightGBM utilizes histogram-based learning to discretize the feature values by binning them, which greatly reduces the memory amount needed for storing feature values and speeds up the learning process, especially when dealing with high-dimensional data.

The combination of the previously mentioned two innovations leaf-wise growth and histogram-based learning enables LightGBM to be faster and more scalable than traditional gradient boosting approaches, especially with large datasets.

#### Performance Evaluation

The LightGBM model was evaluated against a [specific dataset or task], with the following performance metric established:

- **Accuracy:** 58.79%

In this case, the model demonstrated moderate performance, achieving an accuracy of approximately 58.79%. Accuracy, as a primary metric, is calculated as the ratio of correctly predicted instances to the total number of instances. While the model's accuracy is reasonable, the value suggests that there may still be room for improvement, particularly when compared to other models tested in this study.

#### 4.6 Support Vector Classifier

Support Vector Machines (SVMs) are a popular supervised machine learning algorithm used in classification tasks due to their performance in high dimensional spaces and ability to model non-linear relationships in the data. SVMs find a hyper plane or decision boundary that best separates corresponding data points of different classes. This is accomplished by maximizing the margin or distance between the closest points of each class (the support vectors).

The versatility of SVMs is partly due to their ability to use different kernel functions to map data into higher dimensional spaces to enable non-linear decision boundaries. The different types of kernels largely include the following:

- **Linear Kernel:** A linear kernel is used for linearly separable data, where the decision boundary is simply a straight line (or hyperplane in more than two dimensions).
- **Polynomial Kernel:** A polynomial kernel makes it possible to have curved decision boundaries (curved hyperplanes) because it uses a polynomial function to map input features into higher dimensions.
- **Radial Basis Function (RBF) Kernel:** A non-linear kernel can establish some distance metric for computing the similarity of points and essentially create circular (or even spherical) decision boundaries to apply more complex relationships between the classes.

In this project, various SVM models are tested with alternative kernels (linear kernel, polynomial kernel, and RBF kernel) to determine the models' performance on a particular classification task, and then pushing the SVM models

performance by introducing bagging potentially would help reduce variance and improve generalization.

### Performance of Various SVC Models

Let's take a closer look at how different SVC models performed based on their accuracy, which tells us how many predictions were spot on in the dataset:

#### 1. Linear Kernel SVC:

- Accuracy: 66.11%

The linear kernel shines when dealing with data that can be separated in a straight line. With an accuracy of 66.11%, this model did reasonably well, hinting that while there are some linear relationships in the data, there might also be more complex patterns lurking that a simple linear boundary just can't capture.

#### 2. Bagging Classifier with Linear Kernel:

- Accuracy: 65.73%

When we applied bagging to the linear kernel, the performance dipped a bit to 65.73%. Normally, bagging is great for cutting down on model variance by combining predictions from various models trained on different data slices, but here, it didn't really boost the performance compared to the standalone linear kernel SVC.

#### 3. Polynomial Kernel SVC:

- Accuracy: 65.13%

The polynomial kernel offers more flexibility than its linear counterpart, allowing for curved decision boundaries. However, this model's accuracy was a tad lower at 65.13%, suggesting that the data might not have had enough complexity to warrant the extra flexibility that a polynomial kernel brings to the table.

#### 4. Radial Basis Function (RBF) Kernel SVC:

- Accuracy: 66.36%

Known for its knack for capturing non-linear relationships, the RBF kernel achieved an accuracy of 66.36%. This result points to the presence of more intricate, non-linear patterns in the dataset that the RBF kernel could model effectively, leading to a slight edge in accuracy over the linear kernel.

#### 5. Bagging Classifier with RBF Kernel:

- Accuracy: 66.43%

The bagging classifier paired with the RBF kernel hit the highest accuracy at 66.43%. Bagging helped to slightly lower the model's variance while still reaping the benefits of the RBF kernel's ability to handle non-linearities. This small boost in accuracy suggests that bagging improved the model's stability, enhancing its ability to generalize to unseen data.

## 5 Discussion

The comparative analysis of different machine learning models for predicting restaurant ratings reveals some key insights. Among all the models we tested, ensemble methods especially XGBoost and Random Forest with hyperparameter tuning really stood out, outperforming simpler models like

Linear Regression and Support Vector Machines when it came to predictive accuracy. This really highlights how effective tree-based models can be at identifying complex patterns and interactions in customer reviews.

Interestingly, while the sentiment features from TextBlob showed a moderate correlation with actual ratings (with Pearson and Spearman values of 0.52 and 0.49, respectively), they became even more useful when paired with TF-IDF representations. This suggests that basic sentiment analysis alone doesn't quite capture the full depth of customer feedback, but when it's integrated into a broader textual framework, it significantly boosts predictive power.

The performance gains achieved through techniques like SMOTE and class weighting in XGBoost underscore the importance of addressing class imbalance, which is often seen in real-world rating distributions. Plus, models like LightGBM and Extra Trees Regressor delivered competitive results, reinforcing the idea that lightweight or randomized tree-based models can be both effective and computationally efficient alternatives.

However, there are still some limitations to consider. For one, relying on preprocessed review text and handcrafted features might restrict the model's ability to generalize. Additionally, while models like XGBoost provided better accuracy, they often came with higher computational costs. Also, we didn't explore deep learning-based natural language processing models (like BERT), which could potentially enhance our ability to capture semantic context even further.

Future research could build on this work by bringing in transformer-based models, trying out more NLP feature engineering techniques, and adding in user and restaurant metadata like location, cuisine, and pricing to create a more comprehensive prediction. Exploring real-time deployment scenarios, conducting interpretability studies, and incorporating user feedback could also significantly improve how this research is applied in real-world situations.

## References

- [1] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python*. O'Reilly Media Inc., 2009. (*TextBlob & NLP preprocessing foundation*)
- [2] J. Ramos, "Using TF-IDF to Determine Word Relevance in Document Queries," in *Proceedings of the First Instructional Conference on Machine Learning*, 2003. (*TF-IDF technique*)
- [3] P. Domingos, "A few useful things to know about machine learning," *Communications of the ACM*, vol. 55, no. 10, pp. 78–87, 2012.
- [4] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001. (*Random Forest algorithm*)
- [5] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proc. 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2016, pp. 785–794.
- [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002. (*SMOTE implementation for class balancing*)
- [7] G. Ke et al., "LightGBM: A Highly Efficient Gradient Boosting Decision Tree," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017. (*LightGBM*)
- [8] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995. (*SVM foundation for SVC models*)
- [9] L. Breiman, "Bagging Predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996. (*Bagging ensemble method*)

- [10] A. Geron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 2nd ed., O'Reilly Media, 2019.  
*(Covers all models like Linear Regression, Grid Search, SVM, Random Forest, etc.)*
- [11] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.  
*(scikit-learn implementations)*
- [12] B. Liu, "Sentiment Analysis and Opinion Mining," *Synthesis Lectures on Human Language Technologies*, vol. 5, no. 1, pp. 1–167, 2012.  
*(Core sentiment analysis theory relevant to TextBlob)*
- [13] P. Turney, "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews," in *Proc. ACL*, 2002.  
*(Review-based rating prediction)*
- [14] D. Jurafsky and J. H. Martin, *Speech and Language Processing*, 3rd ed., Prentice Hall, 2023.  
*(Advanced NLP concepts)*
- [15] A. Pak and P. Paroubek, "Twitter as a Corpus for Sentiment Analysis and Opinion Mining," in *Proc. LREC*, 2010.  
*(Relevant corpus-based sentiment work)*